

Finding mesoscopic communities in sparse networks

I. Ispolatov[‡], I. Mazo, A. Yuryev

Ariadne Genomics Inc., 9700 Great Seneca Highway, Suite 113, Rockville, Maryland 20850, USA

E-mail: iispolat@lauca.usach.cl

Abstract. We suggest a fast method to find possibly overlapping network communities of a desired size and link density. Our method is a natural generalization of the finite- T superparamagnetic Potts clustering introduced by Blatt, Wiseman, and Domany (Phys. Rev. Lett. **76**, 3251 (1996) and the recently suggested by Reichard and Bornholdt (Phys. Rev. Lett. **93**, 21870 (2004)) annealing of Potts model with global antiferromagnetic term. Similarly to both preceding works, the proposed generalization is based on ordering of ferromagnetic Potts model; the novelty of the proposed approach lies in the adjustable dependence of the antiferromagnetic term on the population of each Potts state, which interpolates between the two previously considered cases. This adjustability allows to empirically tune the algorithm to detect the maximum number of communities of the given size and link density. We illustrate the method by detecting protein complexes in high-throughput protein binding networks.

[‡] Permanent address: Departamento de Física, Universidad de Santiago de Chile, Casilla 302, Correo 2, Santiago, Chile.

1. Introduction

A number of methods have been developed to find clusters, or densely linked communities, in networks. To mention a few, there are clustering algorithms based on link betweenness, number of in-cluster links, random walks, spectrum of connectivity matrix (see a review [1] and [2, 3]), and ordering of spin models [4, 5, 6]. Yet often a need arises to go beyond the existing clustering algorithms as new kinds of communities and networks are analyzed.

Our initial goal was to find protein complexes and functional modules in protein-protein binding networks. Proteins in a complex link together to simultaneously perform a certain function, while members of a functional module sequentially participate in the same cellular process [5]. Both types of clusters usually consist of 10-40 proteins that are stronger linked with each other than with the the rest of the network. Since certain proteins are known to perform functions ubiquitous to several modules, network communities may overlap. We consider protein-protein binding networks of baker yeast and fruit fly, each consisting of $\sim 10^3$ vertices and $\sim 10^3 - 10^4$ links [7, 8, 9]. These networks are composed from the data obtained in Yeast 2-Hybrid (Y2H) high-throughput experiments. Such networks are known to be rather noisy and incomplete, that is, to contain a number of links that do not occur naturally and to miss a noticeable fraction of existing links. Thus it is hard to estimate the precise number of links and nodes that comprise a given protein cluster in such a dataset. Protein binding networks are sparse, so that a probability for an arbitrary pair of nodes to be linked is $\sim 10^{-3}$. While it is assumed that the link density inside a cluster is higher than the average, the precise magnitude of the link density contrast is unknown. Overall, the link density contrast in these networks is relatively low: The largest completely connected subgraph, or clique, contains only four and five vertices in yeast and fly networks, correspondingly. In addition, since many proteins function on their own, there are parts of the network that do not belong to any cluster at all.

To summarize, we looked for an *a priori* unknown number of possibly overlapping mesoscopic clusters in a sparse network with a low link density contrast. Unfortunately, we were unable to detect a sufficient number of such clusters using any of the existing algorithms. Crucial limitations of many of the available network clustering methods are discussed in [6]. For example, for our purposes we ruled out the Q-optimization algorithm by Newman [2] as in its earliest steps it connects all the vertices with a single neighbor (leaves) to their neighbors, thus making it impossible to select only densely linked clusters such as cliques. Similarly, the clustering algorithm based on consecutive cutting the links with the highest betweenness [3] produces the leafy branches as the links leading to leaves have the lowest betweenness and are the last to be cut. A finite-temperature ordering of Potts model used in [5] to detect protein communities yields in our case only very large (≈ 500 vertices) cluster. The main reason for this failure of the finite-temperature Potts model clustering is a difference in the networks: In addition to the links from Y2H experiments, the network analyzed in [5] contained the data

obtained using other methods such as mass spectroscopy, where protein complexes are often recorded as cliques. A clustering based on an annealing in ferromagnetic Potts model with global antiferromagnetic term [6] performs somewhat better; yet it still did not allow us to find the expected number of mesoscopic communities. However, a generalization of the last two approaches enabled us to detect a large number of candidates for protein complexes and modules of the desired size. In the following section we discuss the methods developed in [4, 5, 6] in more detail and introduce our clustering algorithm. In section III we discuss the implementation of the algorithm, averaging, which is used to check robustness of the found complexes, and present examples. A discussion and a brief summary concludes the paper.

2. Ordering of Potts model on a network

First consider a q -state ferromagnetic Potts model on a network. Each vertex is assigned a state σ (often called a spin) that may have any integer value between one and q . The energy of the system is equal to the number of links that connect pairs of vertices in the same state, so that the Hamiltonian reads

$$H = - \sum_{\{i,j\} \in E} \delta_{\sigma_i, \sigma_j}, \quad (1)$$

where sum runs over all edges and the coupling constant is set equal to one. Evidently, in the ground state all connected vertices are in the same Potts state. Equilibration at a low but finite temperature T results in a mosaic of sets of the same-state vertices, interpreted as network communities [4, 5]. Usually performed in the Canonical ensemble, such finite-temperature equilibration minimizes the free energy $F = H - TS$. The entropy S can be qualitatively approximated by its mean-field form (see, for example, [10]),

$$S_{MF} = N \ln N - \sum_{s=1}^q n_s \ln n_s, \quad N = \sum_{s=1}^q n_s, \quad (2)$$

Here n_s is the number of vertices in state s and N is the total number of vertices in the network. This approximation sets an upper limit on the actual entropy of the network Potts model. Yet it illustrates the process of equilibration as a competition between the energy term H , that favors condensation of all spins into a single state ($n_i = N$, $n_j = 0$, $j \neq i$), and an entropic term $T \sum_{s=1}^q n_s \ln n_s$, that favors a completely disordered configuration ($n_s = N/q$, $s = 1, \dots, q$). A similar competition between ordering and disordering trends defines the structure of the ground state of the Potts model with a global antiferromagnetic term suggested in Ref. [6],

$$H' = - \sum_{\{i,j\} \in E} \delta_{\sigma_i, \sigma_j} + \gamma \sum_{s=1}^q \frac{n_s(n_s - 1)}{2}. \quad (3)$$

where γ is an antiferromagnetic coupling constant. To generalize, the ordering in both the finite-temperature Potts model and zero-temperature model (3) corresponds to

minimization of the expression

$$\tilde{H} = - \sum_{\{i,j\} \in E} \delta_{\sigma_i, \sigma_j} + \sum_{s=1}^q n_s f(n_s), \quad (4)$$

with

$$f(x) = \begin{cases} T \ln(x) & \text{in finite-T Potts model} \\ \gamma x/2 & \text{in the model [4].} \end{cases} \quad (5)$$

Terms that depend only on N are left out. The role of the temperature in these two cases is somewhat different: In the former case the temperature is used as an effective disordering (antiferromagnetic) parameter, while in the later case it is a mean to anneal the system into a sufficiently low-energy configuration. It seems natural to interpret two forms of $f(x)$ in (5) as two particular cases of some general antiferromagnetic penalty function with more than one parameter. Furthermore, the existence of only single adjustable parameter in both cases (5) often does not allow to control the properties such as size and link density of the clusters. As we observed, the Potts model [4, 5] on Y2H protein networks at a certain temperature exhibits a sharp transition from a completely disordered state to a state consisting of a single large (containing $\sim 10\%$ or more of all vertices) ordered component and disordered rest of the system. A possible interpretation of such a large-scale ordering is that the dependence of the disordering (entropy) term on the number of each state spins is weak (logarithmic) and the large increase in cluster size does not carry a sufficient free energy penalty. Indeed, the modified Potts model (3), where the dependence of the anti-clustering term on the number of spins in each state is stronger (linear), yielded several smaller clusters. Evidently, the form of $f(x)$ defines the sizes of ordered clusters: The faster $f(x)$ increases with x , the stronger large clusters are suppressed. In order to overcome the limitations of the existing Potts model clustering methods, it appears natural to go beyond two particular forms of f (5). We consider the generalized Potts Hamiltonian (4) where the global antiferromagnetic term that has two adjustable parameters,

$$f(x) = \gamma x^\alpha, \quad \alpha > 0. \quad (6)$$

The clustering methods of [4, 5] and [6] correspond to $\alpha \rightarrow +0$ and $\alpha = 1$ cases, respectively. In a smaller α case larger communities are produced. while a larger α results in a higher number of smaller clusters. In either case γ should be sufficiently small to observe any ordering at all.

To illustrate the clustering, consider the evolution of a single ordered mesoscopic community of n_1 vertices. We assume that the number of Potts states q is much larger than the number of communities and the bulk of the network remains disordered, so that $n_i = (N - n_1)/(q - 1)$, $i = 2, \dots, q$. The antiferromagnetic term for this configuration reads

$$H_{AF} = \gamma \sum_{s=1}^q n_s^{\alpha+1} = \gamma \left[n_1^{\alpha+1} + (q-1) \left(\frac{N - n_1}{q-1} \right)^{\alpha+1} \right]. \quad (7)$$

The community continues to grow while the number of links ΔL brought into the community by Δn_1 added vertices (usually $\Delta n_1 = 1$) exceeds the antiferromagnetic cost of such vertex addition, that is,

$$\frac{\Delta L}{\Delta n_1} \geq \gamma(\alpha + 1) \left[n_1^\alpha - \left(\frac{N - n_1}{q - 1} \right)^\alpha \right]. \quad (8)$$

Adjusting γ and α it is possible to detect communities of a desired size and link density.

Evidently, any finite-temperature system has a certain degree of disorder and consequently, non-zero entropy. However, the contribution of the entropy term to the free energy (1) can be made arbitrary small by annealing the system to the sufficiently low temperature. Comparing the entropy and the antiferromagnetic terms, we estimate the threshold temperature $T^* \approx \gamma n^\alpha / \ln n$. Below T^* the equilibrium ordering of clusters of size n and larger is controlled by competition between only the ferromagnetic and antiferromagnetic couplings, leaving the entropic term irrelevant. This is of course only a qualitative estimate as it is based on a mean-field approximation for the entropy (2).

3. Implementation and Averaging

To make the antiferromagnetic term work, the disordered equilibrium population of a state N/q should be significantly less than a size of the smallest cluster we need to detect. We set $N/5 \leq q \leq N/3$, and experimentally determine the optimal values of γ and α . For the Y2H networks [7, 8, 9] these numbers are $0.002 \leq \gamma \leq 0.02$ and $1 \leq \alpha \leq 2$. In general, we observed that for a typically sparse protein binding network where $2L/N^2 \sim 10^{-3}$, it is convenient to start with $\alpha = 1$ as in [6], adjust $\gamma \sim 10^{-2}$ to produce the reasonable number of clusters, and then fine-tune both α and γ to focus on the desired cluster size and link density. To illustrate this process, cluster abundance vs cluster size plots are presented in Fig. 1 for three sets of (α, γ) . Similarly to [6], the network is initialized with randomly assigned spins and then gradually annealed to $T \ll T^*$. At an annealing step a state of a randomly picked spin is evolved according to the Metropolis rules; each spin is approached at average Cq times with $C \sim 10$. After such isothermal equilibration the temperature is reduced by a small fraction (usually 1-2 %). The algorithm is fairly fast, its performance scales as Nq .

Naturally, each run produces a distinct set of clusters. In some sense, all clusters of the expected size that contain sufficient number of links are good as is, since their high link density make them equally good candidates for protein complexes. However, certain communities are reproduced practically in all runs, while the others are not so robust. Such lack of robustness often has the following explanation: There may exist a set of vertices that contribute similar numbers of links if brought into a cluster. However, in each run only a fraction of these vertices is included into a cluster due to the rapidly growing antiferromagnetic cost (8). Alternating membership of such vertices in a cluster results in its poor reproducibility.

To study the robustness of clusters more systematically, we average the results of many annealing runs. Along with the averaging methods and cluster merging algorithms

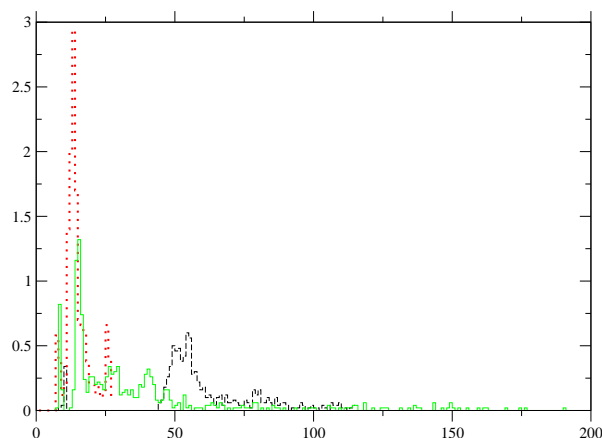


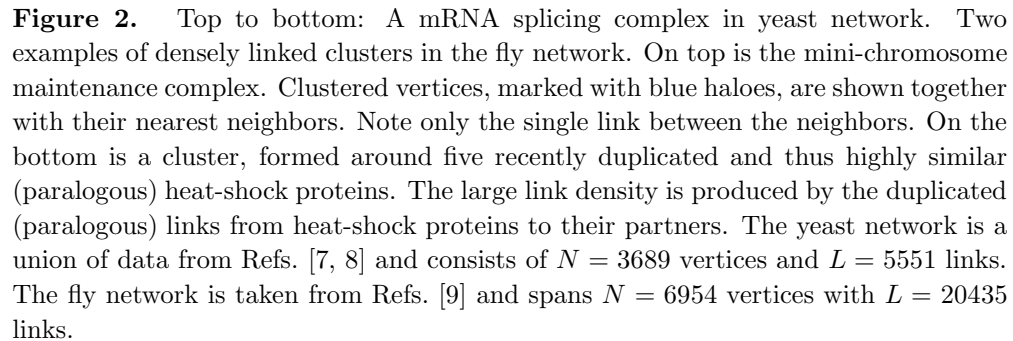
Figure 1. Cluster size histogram for the fly network. The cluster abundance for $\alpha = 1.5$ and $\gamma = 10^{-2}$ (red dotted line) strongly peaks around the desired community size, $n \approx 15$, while the histogram for the same α and smaller $\gamma = 10^{-3}$ (dashed black line) consists of a smaller and broader peak at much larger clusters, $n \approx 50$. While clustering with a smaller antiferromagnetic exponent, $\alpha = 0.5$ and $\gamma = 0.2$ (green solid line) also produces a cluster size distribution with a maximum at a desired cluster size, $n \approx 15$, the number of such clusters is noticeably less than in the $\alpha = 1.5$ and $\gamma = 10^{-2}$ case, and very large (up to $n = 200$) biologically-irrelevant clusters are produced. Only clusters consisting of $n > 8$ vertices and $L > 2n$ links are counted, the results are averaged over 50 equilibration runs.

used in [4, 6], we utilize the following procedure. In each run the “ordered” links that connect the same state vertices are marked. As a result, each link carries an order parameter $\psi \leq 1$ equal to the fraction of runs in which this link was ordered. For a community obtained in a particular annealing run, the averaged over all in-community links order parameter $\bar{\psi}$ characterizes the reproducibility of the community. It was not uncommon to see communities with $\bar{\psi} = 0.5$ and higher.

In each run we were able to detect 5 – 15 (in the baker yeast network) and 15 – 30 (in the fruit fly network) communities of $n > 10$ vertices and $L \geq 2n$ in-community links. Examples of candidates for protein complexes revealed by this algorithm are shown in Fig. 2.

4. Discussion and conclusions

Revealing the intrinsic connection between the finite- T Potts ordering and zero- T Potts clustering with the additional antiferromagnetic coupling [6], we developed a fast method for detecting mesoscopic-size communities in sparse networks. Our method is a natural generalization of the algorithms introduced in [4, 6] and is based on the Potts model with a two-parameter global antiferromagnetic term (4,6). Applying the method to the protein binding networks of the fruit fly and baker yeast, we were able to detect more than a hundred densely interlinked communities that included strong candidates for not yet annotated protein complexes and functional modules.



Acknowledgment

This work was supported by 1 R01 GM068954-01 grant from NIGMS.

References

- [1] M. E. J. Newman, Eur. Phys. J. B **38**, 321-330 (2004).
- [2] M. E. J. Newman, Phys. Rev. E **69**, 066133 (2004).
- [3] M. E. J. Newman, M. Girvan, Phys. Rev. E **69**, 026113 (2004).
- [4] M. Blatt, S. Wiseman, and E. Domany, Phys. Rev. Lett. **76**, 3251 (1996).
- [5] V. Spirin and L.E. Mirny, PNAS **100**, 12123 (2003).
- [6] J. Reichard and S. Bornholdt, Phys. Rev. Lett. **93**, 21870 (2004).
- [7] T. Ito, T. Chiba, R. Ozawa, M. Yoshida, M. Hattori, and Y. Sakaki PNAS **98**, 4569 (2001).
- [8] P. Uetz et al, Nature **403**, 623 (2000).
- [9] L. Giot et al, Science **302**, 1727 (2003).
- [10] F. W. Wu, Rev. Mod. Phys. **54**, 235 (1982).